

# Scientific workflows for process mining: building blocks, scenarios, and implementation

Alfredo Bolt<sup>1</sup> · Massimiliano de Leoni<sup>1</sup> · Wil M. P. van der Aalst<sup>1</sup>

Published online: 14 September 2015

© The Author(s) 2015. This article is published with open access at Springerlink.com

**Abstract** Over the past decade process mining has emerged as a new analytical discipline able to answer a variety of questions based on event data. Event logs have a very particular structure; events have timestamps, refer to activities and resources, and need to be correlated to form process instances. Process mining results tend to be very different from classical data mining results, e.g., process discovery may yield end-to-end process models capturing different perspectives rather than decision trees or frequent patterns. A process-mining tool like *ProM* provides hundreds of different process mining techniques ranging from discovery and conformance checking to filtering and prediction. Typically, a combination of techniques is needed and, for every step, there are different techniques that may be very sensitive to parameter settings. Moreover, event logs may be huge and may need to be decomposed and distributed for analysis. These aspects make it very cumbersome to analyze event logs manually. *Process mining should be repeatable and automated.* Therefore, we propose a framework to support the analysis of process mining workflows. Existing scientific workflow systems and data mining tools are not tailored towards process mining and the artifacts used for analysis (process models and event logs). This paper structures the *basic building blocks* needed for process mining and describes various analysis scenarios. Based on these requirements we implemented *RapidProM*, a tool supporting scientific workflows for process mining. Examples illustrating the different scenarios are provided to show the feasibility of the approach.

**Keywords** Scientific workflows · Process mining · Large scale process analysis · RapidProM

## 1 Introduction

Scientific Workflow Management (SWFM) systems help users to design, compose, execute, archive, and share workflows that represent some type of analysis or experiment. Scientific workflows are often represented as directed graphs where the nodes represent “work” and the edges represent paths along which data and results can flow between nodes. Next to “classical” SWFM systems such as Taverna [23], Kepler [33], Galaxy [20], ClowdFlows [27], and jABC [40], one can also see the uptake of integrated environments for data mining, predictive analytics, business analytics, machine learning, text mining, reporting, etc. Notable examples are RapidMiner [22] and KNIME [4]. These can be viewed as SWFM systems tailored towards the needs of data scientists.

Traditional data-driven analysis techniques do not consider end-to-end processes. People are process models by hand [e.g., Petri nets, UML activity diagrams, or *Business Process Modeling Notation* (BPMN) models], but this modeled behavior is seldom aligned with real-life event data. *Process mining* aims to bridge this gap by connecting end-to-end process models to the raw events that have been recorded.

Process-mining techniques enable the analysis of a wide variety of processes using event data. For example, event logs can be used to automatically learn a process model (e.g., a Petri net or BPMN model). Next to the automated discovery of the real underlying process, there are process-mining techniques to analyze bottlenecks, to uncover hidden inefficiencies, to check compliance, to explain deviations, to predict performance, and to guide users towards “better”

✉ Wil M. P. van der Aalst  
w.m.p.v.d.aalst@tue.nl

<sup>1</sup> Department of Mathematics and Computer Science,  
Eindhoven University of Technology, Eindhoven,  
The Netherlands

processes. Hundreds of process-mining techniques are available and their value has been proven in many case studies. See for example the twenty *case studies* on the webpage of the IEEE Task Force on Process Mining [24]. The open source process mining framework *ProM* [58] provides hundreds of plug-ins and has been downloaded over 100,000 times. The growing number of commercial *process mining tools* (Disco, Perceptive Process Mining, Celonis Process Mining, QPR ProcessAnalyzer, Software AG/ARIS PPM, Fujitsu Interstage Automated Process Discovery, etc.) further illustrates the uptake of process mining.

For process mining typically many analysis steps need to be chained together. Existing process mining tools do not support such analysis workflows. As a result, analysis may be tedious and it is easy to make errors. Repeatability and provenance are jeopardized by manually executing more involved process mining workflows.

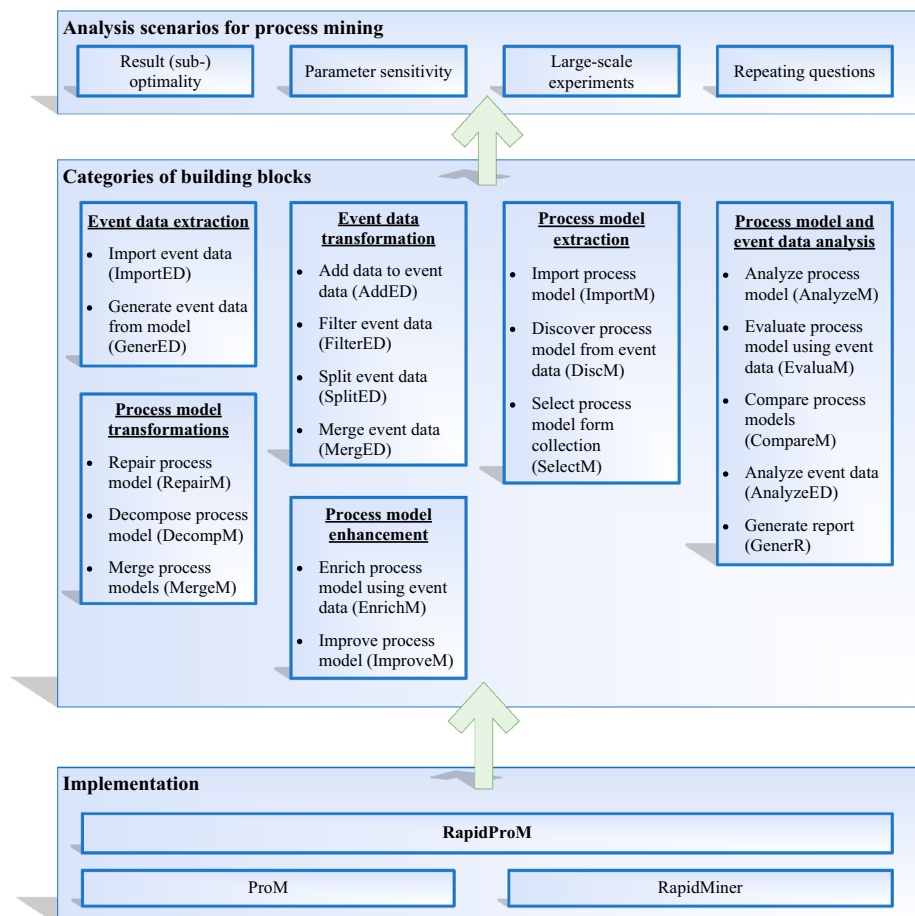
This paper is motivated by the observation that tool support for process mining workflows is missing. None of the process mining tools (ProM, Disco, Perceptive, Celonis, QPR, etc.) provides a facility to design and execute analysis workflows. None of the scientific workflow management systems including analytics suites like RapidMiner and KNIME support process mining. Yet, process models and event logs

are very different from the artifacts typically considered. Therefore, we propose the *framework to support process mining workflows* depicted in Fig. 1.

This paper considers four *analysis scenarios* where process mining workflows are essential:

- *Result (sub-)optimality* Often different process mining techniques can be applied and a priori it is not clear which one is most suitable. By modeling the analysis workflow, one can just perform all candidate techniques on the data, evaluate the different analysis results, and pick the result with the highest quality (e.g., the process model best describing the observed behavior).
- *Parameter sensitivity* Different parameter settings and alternative ways of filtering can have unexpected effects. Therefore, it is important to see how sensitive the results are (e.g., leaving out some data or changing a parameter setting a bit should not change the results dramatically). It is important to not simply show the analysis result without having some confidence indications.
- *Large-scale experiments* Each year new process mining techniques become available and larger data sets need to be tackled. For example, novel discovery techniques need to be evaluated through massive testing and larger event

**Fig. 1** Overview of the framework to support process mining workflows



logs need to be decomposed to make analysis feasible. Without automated workflow support, these experiments are tedious, error-prone, and time consuming.

- *Repeating questions* It is important to lower the threshold for process mining to let non-expert users approach it. Questions are often repetitive, e.g., the same analysis is done for a different period or a different group of cases. Process mining workflows facilitate recurring forms of analysis.

As shown in Fig. 1 these scenarios build on process mining *building blocks* grouped into six categories:

- *Event data extraction* Building blocks to extract data from systems or to create synthetic data.
- *Event data transformation* Building blocks to pre-process data (e.g., splitting, merging, filtering, and enriching) before analysis.
- *Process model extraction* Building blocks to obtain process models, e.g., through discovery or selection.
- *Process model and event analysis* Building blocks to evaluate event logs and models, e.g., to check the internal consistency or to check conformance with respect to an event log.
- *Process model transformations* Building blocks to repair, merge or decompose process models.
- *Process model enhancement* Building blocks to enrich event logs with additional perspectives or to suggest process improvements.

Building blocks can be chained together to support specific analysis scenarios. The suggested approach has been implemented thereby building on the process mining framework *ProM* and the workflow and data mining capabilities of *RapidMiner*. The resulting tool is called, *RapidProM*, which supports process mining workflows. *ProM* was selected because it is open source and there is no other tool that supports as many process mining building blocks. *RapidMiner* was selected because it allows for extensions that can be offered through a marketplace. *RapidProM* is also offered as such an extension and the infrastructure allows us to mix process mining with traditional data mining approaches, text mining, reporting, and machine learning. Overall, *RapidProM* offers comprehensive support for any type of analysis involving event data and processes.

The remainder of this paper is organized as follows: Section 2 discusses related work and positions our framework. An initial set of process-mining building blocks is described in Sect. 3. These building blocks support the four analysis scenarios described in Sect. 4. The *RapidProM* implementation is presented in Sect. 5. Section 6 evaluates the approach by showing concrete examples. Finally, Sect. 7 concludes the paper.

## 2 Related work

Over the past decade, *process mining* has emerged as a new scientific discipline at the interface between process models and event data [45]. Conventional Business Process Management (BPM) [46, 63] and Workflow Management (WfM) [31, 51] approaches and tools are mostly model-driven with little consideration for event data. Data Mining (DM) [21], Business Intelligence (BI), and Machine Learning (ML) [35] focus on data without considering end-to-end process models. Process mining aims to bridge the gap between BPM and WfM on the one hand and DM, BI, and ML on the other hand. A wealth of *process discovery* [29, 53, 62] and *conformance checking* [1, 2, 48] techniques has become available. For example, the process mining framework *ProM* [58] provides hundreds of plug-ins supporting different types of process mining (<http://www.processmining.org>).

This paper takes a *different* perspective on the gap between analytics and BPM/WfM. We propose to use workflow technology for process mining rather than the other way around. To this end, we focus on particular kinds of scientific workflows composed of process mining operators.

Differences between scientific and business workflows have been discussed in several papers [3]. Despite unification attempts (e.g., [38]) both domains have remained quite disparate due to differences in functional requirements, selected priorities, and disjoint communities.

Obviously, the work reported in this paper is closer to scientific workflows than business workflows (i.e., traditional BPM/WfM from the business domain). Numerous *Scientific Workflow Management (SWFM)* systems have been developed. Examples include Taverna [23], Kepler [33], Galaxy [20], CloudFlows [27], jABC [40], Vistrails, Pegasus, Swift, e-BioFlow, VIEW, and many others. Some of the SWFM systems (e.g., Kepler and Galaxy) also provide repositories of models. The website <http://www.myExperiment.org> lists over 3500 workflows shared by its members [19]. The diversity of the different approaches illustrates that the field is evolving in many different ways. We refer to the book [41] for an extensive introduction to SWFM.

An approach to mine process models for scientific workflows (including data and control dependencies) was presented in [65]. This approach uses “process mining for scientific workflows” rather than applying scientific workflow technology to process mining. The results in [65] can be used to recommend scientific workflow compositions based on actual usage. To our knowledge, *RapidProM* is the only approach supporting “scientific workflows for process mining”. The demo paper [34] reported on the first implementation. In the meantime, *RapidProM* has been refactored based on various practical experiences.

There are many approaches that aim to analyze repositories of scientific workflows. In [64], the authors provide an extensible process library for analyzing jABC workflows empirically. In [14] graph clustering is used to discover sub-workflows from a repository of workflows. Other analysis approaches include [16, 32], and [61].

Scientific workflows have been developed and adopted in various disciplines, including physics, astronomy, bioinformatics, neuroscience, earth science, economics, health, and social sciences. Various collections of reusable workflows have been proposed for all of these disciplines. For example, in [42] the authors describe workflows for quantitative data analysis in the social sciences.

The boundary between data analytics tools and scientific workflow management systems is not well-defined. Tools like RapidMiner [22] and KNIME [4] provide graphical workflow modeling and execution capabilities. Even the scripting in R [25] can be viewed as primitive workflow support. In this paper we build on RapidMiner as it allows us to mix process mining with data mining and other types of analytics. Earlier we developed extensions of ProM for chaining process mining plug-ins together, but these were merely prototypes. We also realized a prototype using an integration between KNIME and ProM. However, for reasons of usability, we opted for RapidMiner as a platform to expose process mining capabilities.

### 3 Definition of the process-mining building blocks

To create scientific workflows for process mining we need to define the building blocks, which are, then, connected with each other to create meaningful analysis scenarios. This section discusses a taxonomy and a repertoire of such building blocks inspired by the so-called “BPM use cases”, which were presented in [46]. The Process-Mining Building Blocks (PMBB) are characterized by two main aspects. First, they are *abstract* as they are not linked to any specific technique or algorithm. Second, they represent logical units of work, i.e., they cannot be conceptually split while maintaining their generality. This does not imply that concrete techniques that implement process-mining building blocks cannot be composed by micro-steps, according to the implementation and design that was used.

The process-mining building blocks can be chained, thus producing process-mining scientific workflows to answer a variety of process-mining questions.

Each process-mining building block takes a number of inputs and produces certain outputs. The input elements represent the set (or sets) of abstract objects required to perform the operation. The process-mining building block component represents the logical unit of work needed to process the inputs and produce the outputs. Inputs and outputs are

indicated through circles, whereas a process-mining building block is represented by a rectangle. Arcs are used to connect the blocks to the inputs and outputs. A generic example of a building block interacting with inputs and outputs is shown in Fig. 2.

Two process-mining building blocks *a* and *b* are chained if one or more outputs of *a* are used as an inputs in *b*. As mentioned, inputs and outputs are depicted by circles. The letter inside a circle specifies the type of the input or output. The following types of inputs and outputs are considered in this paper:

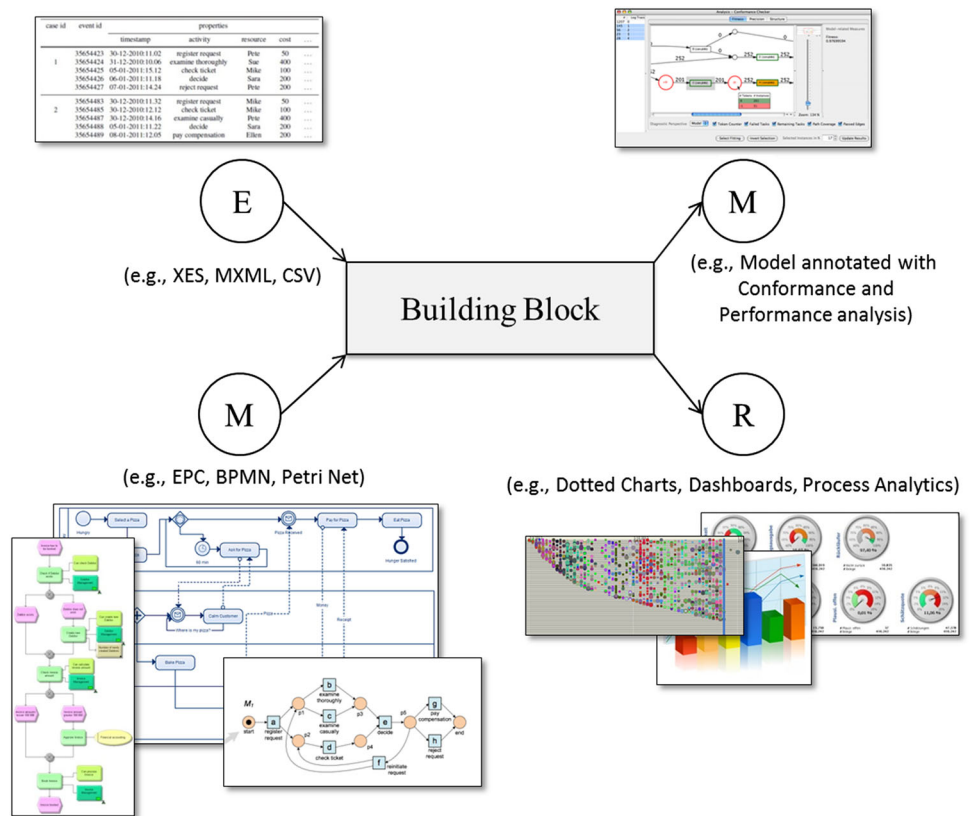
- *Process models*, which are a representation of the behavior of a process, are represented by letter “*M*”. Here we abstract from the notation used, e.g., Petri nets, Heuristics nest, BPMN models are concrete representation languages.
- *Event data sets*, which contain the recording of the execution of process instances within the information system(s), regardless of the format. They are represented by letter “*E*”. XES is currently the de-facto standard format to store events.<sup>1</sup>
- *Information systems*, which supports the performance of processes at runtime. They are represented by the label “*S*”. Information systems may generate events used for analysis and process mining results (e.g., prediction) may influence the information system.
- *Sets of parameters* to configure the application of process-mining building blocks (e.g., thresholds, weights, ratios, etc.). They are represented by letter “*P*”.
- *Results* that are generated as outputs of a process-mining building blocks. This can be as simple as a number or more complex structures like a detailed report. In principle, the types enumerated above in this list (e.g., process models) can also be results. However, it is worth differentiating those specific types of outputs from results which are not process mining specific (like a bar chart). Results are represented by letter “*R*”.
- *Additional Data Sets* that can be used as input for certain process-mining building blocks. These are represented by the letter “*D*”. Such an additional data set can be used to complement event data with context information (e.g., one can use weather or stock-market data to augment the event log with additional data).

The remainder of this section provides a taxonomy of process-mining building blocks grouped into six different categories. For each category, several building blocks are provided. They were selected because of their usefulness

<sup>1</sup> XES (Extensible Event Stream) is an XML-based standard for event logs <http://www.xes-standard.org>. It provides a standard format for the interchange of event log data between tools and application domains.



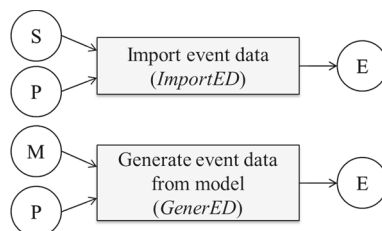
**Fig. 2** Generic example of a building block transforming a process model (M) and event data (E) into process analytics results (R) and an annotated process model



for the definition of many process-mining scientific workflows. The taxonomy is not intended to be exhaustive; there will be new process-mining building blocks as the discipline evolves. Section 5 discusses how these building blocks can be implemented into concrete operators and provides examples of these operators implemented in RapidProM.

### 3.1 Event data extraction

Event data are the *cornerstone* of process mining. In order to be used for analysis, event data has to be extracted and made available. All of the process-mining building blocks of this category can extract event data from different sources. Figure 3 shows some process-mining building blocks that belong to this category.



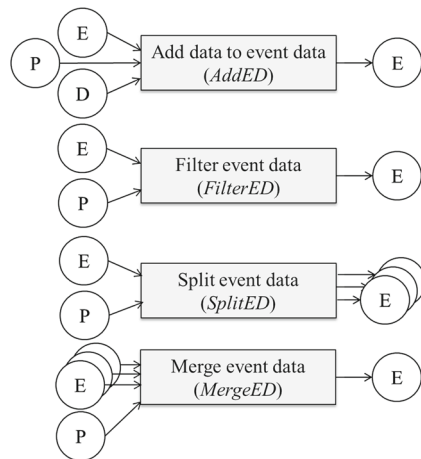
**Fig. 3** Process-mining building blocks related to event data extraction

**Import event data (ImportED)** Information systems store event data in different format and media, from files in a hard drive to databases in the cloud. This building block represents the functionality of extracting event data from any of these sources. Some parameters can be set to drive the event-data extraction. For example, event data can be extracted from files in standard formats, such as XES, or from transactional databases.

**Generate event data from model (GenerED)** In a number of cases, one wants to assess whether a certain technique returns the expected or desired output (i.e., synthetic event data). For this assessment, controlled experiments are necessary where input data are generated in a way that the expected output of the technique is clearly known. Given a process model  $M$ , this building block represents the functionality of generating event data that record the possible execution of instances of  $M$ . This is an important function for, e.g., testing a new discovery technique. Various simulators have been developed to support the generation of event data.

### 3.2 Event data transformation

Sometimes, event data sets are not sufficiently rich to enable certain process-mining analyses. In addition, certain dataset portions should be excluded, because they are irrelevant,



**Fig. 4** Process-mining building blocks related to event data transformations

out of the scope of the analysis or, even, noise. Therefore, a number of event data transformations may be required before doing further analysis. This category comprises the building blocks to provide functionalities to perform the necessary event data transformations. Figure 4 shows the repertoire of process-mining building blocks that belong to this category.

**Add data to event data (AddED)** In order to perform a certain analysis or to improve the results, the event data can be augmented with additional data coming from different sources. For instance, if the process involves citizens, the event data can be augmented with data from the municipality data source. If the level of performance of a process is suspected to be influenced by the weather, event data can incorporate weather data coming from a system storing such a kind of data. If the event data contain a ZIP code, then other data fields such as country or city can be added to the event data from external data sources. This building block represents the functionality of augmenting event data using external data, represented as a generic data set in the figure.

**Filter event data (FilterED)** Several reasons may exist to filter out part of the event data. For instance, the process behavior may exhibit *concept drifts* over time. In those situations, the analysis needs to focus on certain parts of the event data instead of all of it. One could filter the event data and use only those events that occurred, e.g., in year 2015. As a second example, the same process may run at different geographical locations. One may want to restrict the scope of the analysis to a specific location by filtering out the event data referring to different locations. This motivates the importance of being able to filter event data in various ways.

**Split event data (SplitED)** Sometimes, the organization generating the event data is interested in comparing the process'

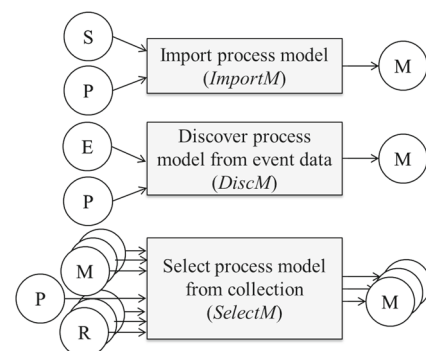
performances for different customers, offices, divisions, involved employees, etc. To perform such comparison, the event data need to be split according to a certain criterion, e.g., according to organizational structures, and the analysis needs to be iterated over each portion of the event data. Finally, the results can be compared to highlight difference. Alternatively, the splitting of the data may be motivated by the size of the data. It may be intractable to analyze all data without decomposition or distribution. Many process-mining techniques are exponential in the number of different activities and linear in the size of the event log. If data are split in a proper way, the results of applying the techniques to the different portions can be fused into a single result. For instance, work [47] discusses how to split event data while preserving the correctness of results. This building block represents the functionality of splitting event data into overlapping or non-overlapping portions.

**Merge event data (MergeED)** This process-mining building block is the inverse of the previous: data sets from different information systems are merged into a single event data set. This process-mining building block can also tackle the typical problems of data fusion, such as redundancy and inconsistency.

### 3.3 Process model extraction

Process mining revolves around process models to represent the behavior of a process. This category is concerned with providing building blocks to mine a process model from event data as well as to select or extract it from a process-model collection. Figure 5 lists a number of process-mining building blocks belonging to this category.

**Import process model (ImportM)** Process models can be stored in some media for later retrieval to conduct some analyses. This building block represents the functionality of loading a process model from some repository.



**Fig. 5** Process-mining building blocks related to process model extraction

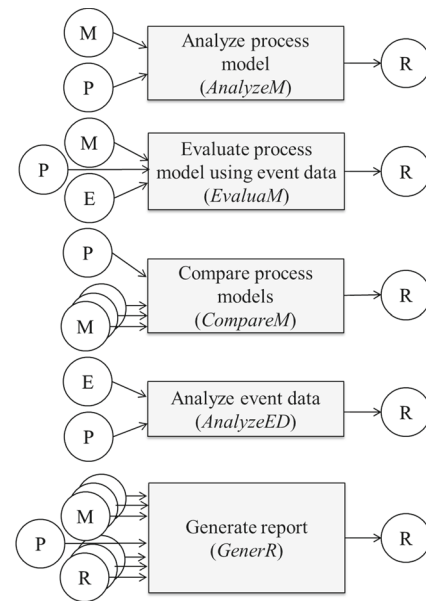
**Discover process model from event data (DiscM)** Process models can be manually *designed* to provide a normative definition for a process. These models are usually intuitive and understandable, but they might not describe accurately what happens in reality. Event data represent the “real behavior” of the process. Discovery techniques can be used to mine a process model on the basis of the behavior observed in the event data (cf. [45]). Here, we stay independent of the specific notations and algorithms. Examples of algorithms are the Alpha Miner [53], the Heuristics Miner [62] or, more recent techniques like the Inductive Miner [29]. This building block represents the functionality of discovering a process model from event data. This block, as many others, can receive a set of parameters as an input to customize the application of the algorithms.

**Select process model from collection (SelectM)** Organizations can be viewed as a collection of processes and resources that are interconnected to form a *process ecosystem*. This collection of processes can be managed and supported by different approaches, such as ARIS [36] or Apromore [28]. To conduct certain analyses, one needs to use some of these models and not the whole collection. In addition, one can give a criterion to retrieve a subset of the collection. This building block represents the functionality of selecting one or more process models from a process-model collection.

### 3.4 Process model and event analysis

Organizations normally use process models for the discussion, configuration, and implementation of processes. In recent years, many process mining techniques are also using process models for analysis. This category groups process-mining building blocks that can analyze process models or event logs and provide analysis results. Figure 6 shows some process-mining building blocks that belong to this category.

**Analyze process model (AnalyzeM)** Process models may contain a number of structural problems. For instance, the model may exhibit undesired deadlocks, activities that are never enabled for execution, variables that are used to drive decisions without previously taking on a value, etc. Several techniques have been designed to verify the soundness of process models against deadlocks and other problems [52]. This building block refers to design-time properties: the process model is analyzed without considering how the process instances are actually being executed. The checking of the conformance of the process model against real event data is covered by the next building block (*Evaluam*). Undesired design-time properties happen for models designed by hand but also for models automatically mined from event data. Indeed, several discovery techniques do not guarantee to mine process models without structural problems. This



**Fig. 6** Process-mining building blocks related to process model and event analysis

building block provides functionalities for analyzing process models and detecting structural problems.

**Evaluate process model using event data (Evaluam)** Besides structural analysis, process models can also be analyzed against event data. Compared with the previous building block (*AnalyzeM*), this block is not concerned with a design-time analysis. Conversely, it makes a-posteriori analysis where the adherence of the process model is checked with respect to the event data, namely how the process has actually been executed. In this way, the expected or normative behavior as represented by the process model is checked against the actual behavior as recorded in the event data. In the literature, this is referred to as *conformance checking* (cf. [45]). This can be used, for example, in fraud or anomaly detection. Replaying event data on process models has many possible uses: aligning observed behavior with modeled behavior is key in many applications. For example, after aligning event data and model, one can use the *time* and *resource* information contained in the log for performance analysis. This can be used for *bottleneck* identification or to gather information for simulation analysis or predictive techniques. This building block represents the functionality of analyzing or evaluating process models using event data.

**Compare process models (CompareM)** Processes are not static as they dynamically evolve and adapt to the business context and requirements. For example, processes can behave differently over different years, or at different locations. Such differences or similarities can be captured through the comparison of the corresponding process models. For example,

the degree of similarity can be calculated. Approaches that explicitly represent configuration or variation points [49] directly benefit from such comparisons. Building block *CompareM* is often used in combination with *SplitED* that splits the event data into sublogs and *DiscM* that discovers a model per sublog.

**Analyze event data (AnalyzeED)** Instead of directly creating a process model from event data, one can also first inspect the data and look at basic statistics. Moreover, it often helps to simply visualize the data. For example, one can create a so-called dotted chart [45] exploiting the temporal dimension of event data. Every event is plotted in a two-dimensional space where one dimension represents the time (absolute or relative) and the other dimension may be based on the case, resource, activity or any other property of the event. The color of the dot can be used as a third dimension. See [26] for other approaches combining visualization with other analytical techniques.

**Generate report (GenerR)** To consolidate process models and other results, one may create a structured report. The goal is not to create new analysis results, but to present the findings in an understandable and predictable manner. Generating standard reports helps to reduce the cognitive load and helps users to focus on the things that matter most.

### 3.5 Process model transformations

Process models can be designed or, alternatively, discovered from event data. Sometimes, these models need to be adjusted for follow-up analyses. This category groups process-mining building blocks that provide functionality to change the structure of a process model. Figure 7 shows some process-mining building blocks that belong to this category.

**Repair process model (RepairM)** Process models may need to be repaired in case of consistency or conformance problems. Repairing can be regarded from two perspectives:

repairing structural problems and repairing behavioral problems. The first case is related to the fact that models can contain undesired design-time properties such as deadlocks and livelocks (see also the *Analyze process model* building block discussed in Sect. 3.4). Repairing involves modifying the model to avoid those properties. Techniques for repairing behavioral problems focus on models that are structurally sound but that allow for undesired behavior or behavior that does not reflect reality. See also the *Evaluate process model using event data* building block discussed in Sect. 3.4, which is concerned with discovering the conformance problems. This building block provides functionality for both types of repair.

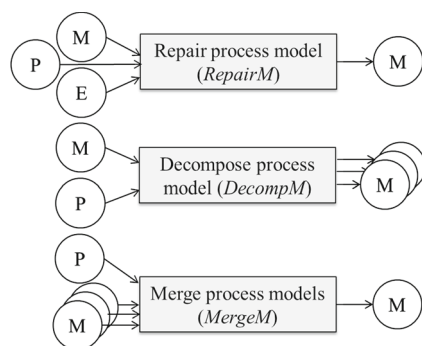
**Decompose process model (DecompM)** Processes running within organizations may be extremely large, in terms of activities, resources, data variables, etc. As mentioned, many techniques are exponential in the number of activities. The computation may be improved by splitting the models into fragments, analogously to what was mentioned for splitting the event log. If the model is split according to certain criteria, the results can be somehow amalgamated and, hence, be meaningful for the entire model seen as a whole. For instance, the work on decomposed conformance checking [47] discusses how to split process model to make process mining possible with models with hundreds of elements (such as activities, resources, data variables), while preserving the correctness of certain results (e.g., the fraction of deviating cases does not change because of decomposition). This block provides functionalities for splitting process models into smaller fragments.

**Merge process models (MergeM)** Process models may also be created from the intersection (i.e., the common behavior) or union of other models. This building block provides functionalities for merging process models into a single process model. When process discovery is decomposed, the resulting models need to be merged into a single model.

### 3.6 Process model enhancement

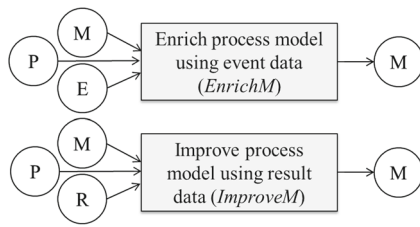
Process models just describing the control-flow are usually not the final result of process mining analysis. Process models can be enriched or improved using additional data to provide better insights into the real process behavior that it represents. This category groups process-mining building blocks that are used to enhance process models. Figure 8 shows a summary of the process-mining building blocks that belong to this category.

**Enrich process model using event data (EnrichM)** The backbone of any process models contains basic structural information relating to control-flow. However, the backbone can



**Fig. 7** Process-mining building blocks related to process model transformations





**Fig. 8** Process-mining building blocks related to process model enhancement

be enriched with additional perspectives derived from event data to obtain better analysis results. For example, event frequency can be annotated in a process model to identify the most common paths followed by process instances. Timing information can also be used to enrich a process model to highlight bottlenecks or long waiting times. This enrichment does not have an effect on the structure of the process model. This building block represents the functionality of enriching process models with additional information contained in event data.

*Improve process model (ImproveM)* Besides being enriched with data, process models can also be improved. For example, performance data can be used to suggest structural modifications to improve the overall process performance. It is possible to automatically improve models using causal dependencies and observed performance. The impact of such modifications could be simulated in “what-if scenarios” using performance data obtained in the previous steps. This building block represents the functionality of improving process models using data from other analysis results.

## 4 Analysis scenarios for process mining

This section identifies generic analysis scenarios that are not domain-specific and, hence, that can be applied to different contexts. The analysis scenarios compose the basic process-mining building blocks and, hence, they remain independent of any specific operationalization of a technique. In fact, as mentioned before, the building blocks may employ different concrete techniques, e.g., there are dozens of process discovery techniques realizing instances of building block *DiscM* (Fig. 5).

As depicted in Fig. 1, we consider four analysis scenarios: (a) *result (sub-)optimality*, (b) *parameter sensitivity*, (c) *large-scale experiments*, and (d) *repeating questions*. These are described in the remainder of this section.

As discussed in this section and validated in Sect. 6, the same results could also be achieved without using scientific workflows. However, the results would require a tedious and error-prone work of repeating the same steps *ad nauseam*.

### 4.1 Result (sub-)optimality

This subsection discusses how process-mining building blocks can be used to mine optimal process models according to some optimality criteria. Often, in process discovery, optimality is difficult (or even impossible) to achieve. Often sub-optimal results are returned and it is no known what is “optimal”.

Consider for example the process discovery task. The quality of a discovered process model is generally defined by four quality metrics [1,2,45,48]:

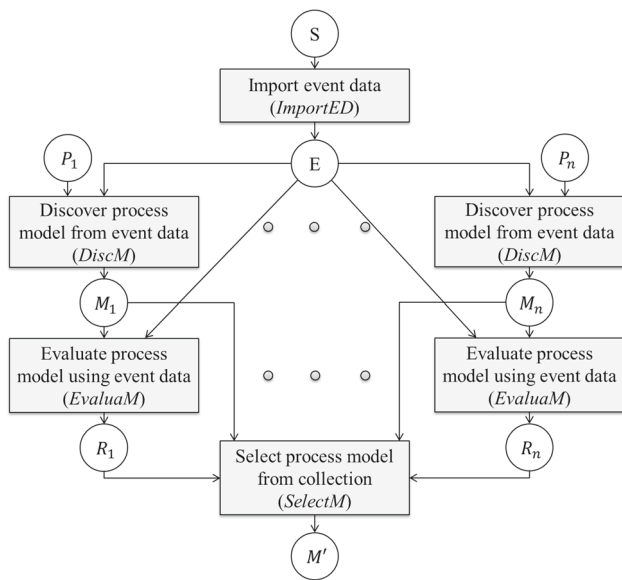
- *Replay fitness* quantifies the ability of the process model to reproduce the execution of process instances as recorded in event data.
- *Simplicity* captures the degree of complexity of a process model, in terms of the numbers of activities, arcs, variables, gateways, etc.
- *Precision* quantifies the degree with which the model allows for too much behavior compared to what was observed in the event data.
- *Generalization* quantifies the degree with which the process model is capable to reproduce behavior that is not observed in the event data but that potentially should be allowed. This is linked to the fact that event data often are incomplete in the sense that only a fraction of the possible behaviors can be observed.

Traditionally, these values are normalized between 0 and 1, where 1 indicates the highest score and 0 the lowest.

The model of the highest value within a collection of (discovered) models is such that it can mediate among those criteria. Often, these criteria are in competing: higher score for one criterion may lower the score of a second criterion. For instance, to have a more precise model, it is necessary to sacrifice the behavior observed in the event data that is less frequent, thus partly hampering the replay-fitness score.

Figure 9 shows a suitable scientific workflow for mining a process model from event data that is sub-optimal with respect to a score defined by specific criteria. The optimization is done by finding the parameters that returns a sub-optimal model.

Event data are loaded from an information system and used  $n$  times as input for a discovery technique using different parameter values. The  $n$  resulting process models are evaluated using the original event data and the model that has the best score is returned. Please note that the result is likely to be sub-optimal:  $n$  arbitrary parameter values are chosen out of a much larger set of possibilities. If  $n$  is sufficiently large, the result is sufficiently close to the optimal. This scientific workflow is still independent of the specific algorithm used for discovery; as such, the parameter settings are also generic.



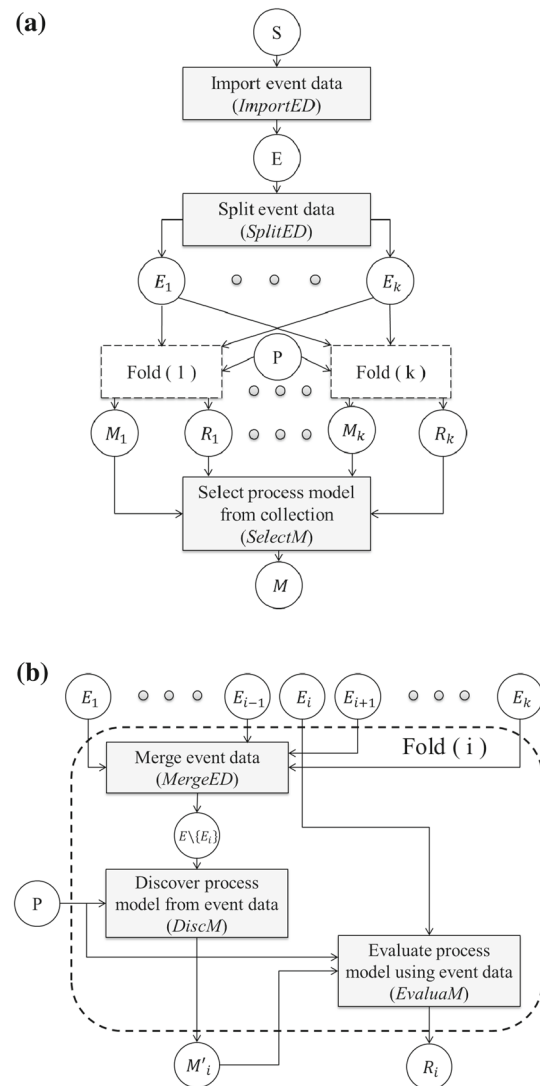
**Fig. 9** Result (sub-)optimality in process model discovery: process-mining scientific workflow for mining an optimal model in terms of a defined scoring criteria

Figure 10a illustrates a scientific workflow that tries to account for generalization. For this purpose, a *k-fold cross validation* approach is used. In this approach, the process instances recorded in the event data are randomly split into *k* folds, through building block *Split event data* (*SplitED*). For each of the *k* times, a different fold is taken aside: the other *k* − 1 folds are used for discovery and the “elected” fold is used for evaluation through conformance checking. This corresponds to block *Fold(i)* with  $1 \leq i \leq k$ . Finally, through the process-mining building block *Select process model from collection* (*SelectM*), the model with the best score is returned as output. Figure 10b enters inside the block *Fold(i)* showing how fold  $E_i$  is used for evaluation and folds  $E_1, \dots, E_{i-1}, E_{i+1}, E_n$  are used for discovery (after being merged).

Scientific workflows can also be hierarchically defined: in turn, the discover process-mining building block (*DiscM*) in Fig. 9 can be an entire scientific sub-workflow. The two scientific workflows shown in Figs. 9 and 10 do not exclude each other. Process-mining building block *Discover process model from event data* (*DiscM*) can be replaced by the entire workflow in Fig. 10a, thus including some generalization aspects in the search for a sub-optimal process model.

## 4.2 Parameter sensitivity

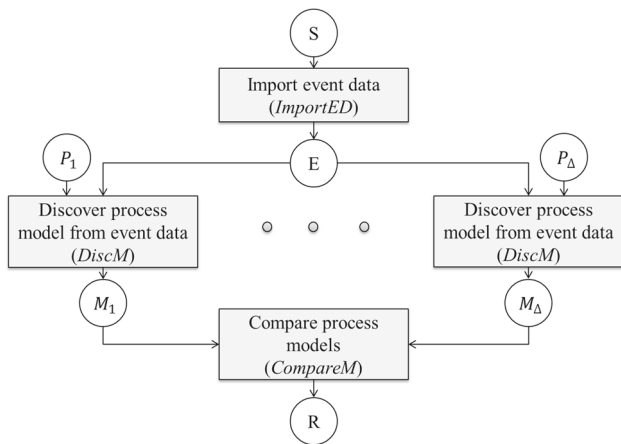
Parameters are used by techniques to customize their behavior, e.g., adapting to the noise level in the event log. These parameters have different ways of affecting the results produced, depending on the specific implementation of the technique or algorithm. Some parameters can have more *rel-*



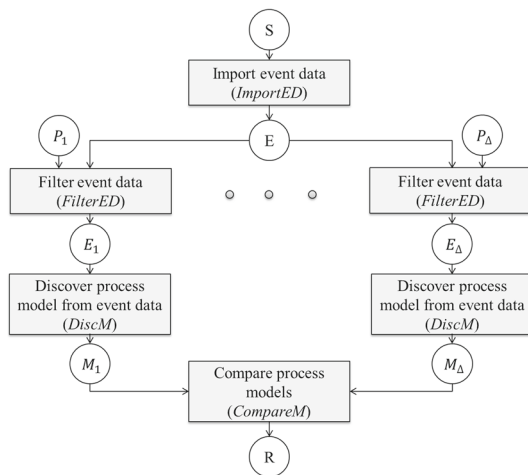
**Fig. 10** Process-mining main scientific workflow based on *k*-fold cross validation. **a** Main workflow. **b** Process-mining sub-workflow for any *Fold(i)*

*evance* than others (i.e., they have a more substantial effect on the results). There are many ways to evaluate the sensitivity of a certain parameter for a given algorithm. Figure 11 shows an example of this analysis scenario. Here the parameter value is varied across the range. For each of the discovered models, a score is computed. The results are finally plotted on a Cartesian coordinate system where the X-axis is associated with the potential parameter’s values and the Y-axis is associated with the model’s score.

Alternatively, the sensitivity analysis can also focus on the filtering part, while keeping the same configuration of parameter(s) for discovery. In other words, we can study how the discovered model is affected by different filtering, namely different values of the parameter(s) that customize the application of filtering.



**Fig. 11** Parameter sensitivity in process discovery techniques: process mining workflow for comparing the effects of different parameter values for a given discovery technique



**Fig. 12** Parameter sensitivity in event data filtering: process-mining scientific workflow for comparing the effect of different event-data filtering configurations on the discovered model

Figure 12 shows an example of this analysis scenario in the process mining domain, using process-mining building block to analyze the differences and similarities of results obtained by discovery techniques from event data that were filtered using different parameter values. In this example, event data are loaded and filtered several times using different parameter settings, producing several filtered event data sets. Each of these filtered event data sets is input for the same discovery technique using the same configuration of parameter(s).

### 4.3 Large-scale experiments

Empirical evaluation is often needed (and certainly recommended) when testing new process mining algorithms. In case of process mining, many experiments need to be conducted to prove that these algorithms or techniques can be

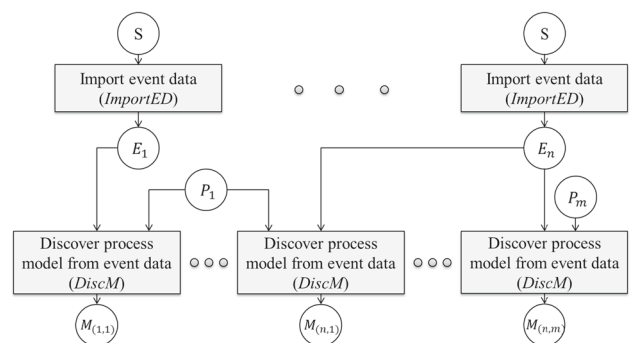
applied in reality and that the results are as expected. This is due to the richness of the domain. Process models can have a wide variety of routing behaviors, timing behavior, and second-order dynamics (e.g., concept drift). Event logs can be large or small and contain infrequent behavior (sometimes called noise) or not. Hence this type of evaluation has to be conducted on a large scale. The execution and evaluation of such large-scale experiment results is a tedious and time-consuming task: it requires intensive human assistance by configuring each experiment's run and waiting for the results at the end of each run.

This can be greatly improved by using process mining workflows, as only one initial configuration is required. There are many examples for this analysis scenario within the process mining domain. Two of them are presented next.

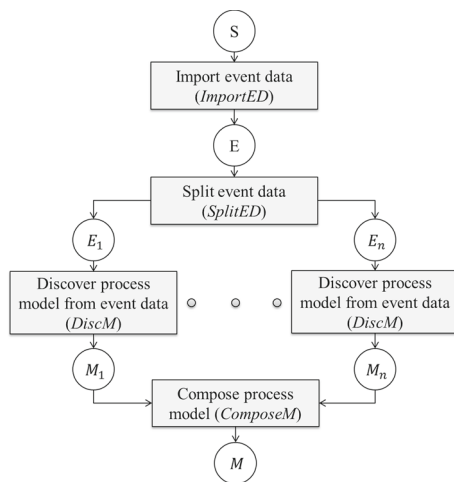
#### 4.3.1 Assessment of discovery techniques through massive testing

When developing new process discovery techniques, several experiments have to be conducted to test the *robustness* of the approach. As mentioned, many discovery techniques use parameters that can affect the result produced. It is extremely time-consuming and error prone to assess the discovery techniques using several different combinations of parameter values and, at the same time, testing on a dozen of different event-data sets.

Figure 13 shows the result of a large-scale experiment using  $n$  event data sets and  $m$  different parameter settings that produces  $n \times m$  resulting process models. In this example, the same discovery technique with different parameters is used. However, one can consider the discovery algorithm to employ as an additional parameter. Therefore, the  $m$  different parameter settings can indicate  $m$  different discovery algorithms. After mining  $n \times m$  models, the best model is considered.



**Fig. 13** Exhaustive testing of a discovery technique: large-scale experiments using different types of event data and parameter combinations are needed to evaluate a discovery technique



**Fig. 14** Decomposed process discovery: a generic example using event data splitting, model composition and a specified discovery technique

### 4.3.2 Decomposed process discovery

Existing process mining techniques are often unable to handle “big event data” adequately. Decomposed process mining aims to solve this problem by decomposing the process mining problem into many smaller problems, which can be solved in less time and using less resources.

In decomposed process discovery, large event data sets are decomposed in sublogs, each of which refers to a subset of the process’ activities. Once an appropriate decomposition is performed, the discovery can be applied to each cluster. This results in as many process models as the number of clusters; these models are finally merged to obtain a single process model. See for example the decomposed process mining technique described in [59] which presents an approach that clusters the event data, applies discovery techniques to each cluster, and merges the process models.

Figure 14 shows a process-mining workflow that splits the event data into  $n$  subsets, then uses a discovery algorithm to discover models for each of these subsets, and finally merges them into a single process model.

## 4.4 Repeating questions

Whereas the previous scenarios are aimed at (data) scientists, process mining workflows can also be used to lower the threshold for process mining. After the process mining workflow has been created and tested, the same analysis can be repeated easily using different subsets of data and different time-periods. Without workflow support this implies repeating the analysis steps manually or using *hardcoded* scripts that perform them over some input data. The use of scientific workflows is clearly beneficial: the same workflow can be

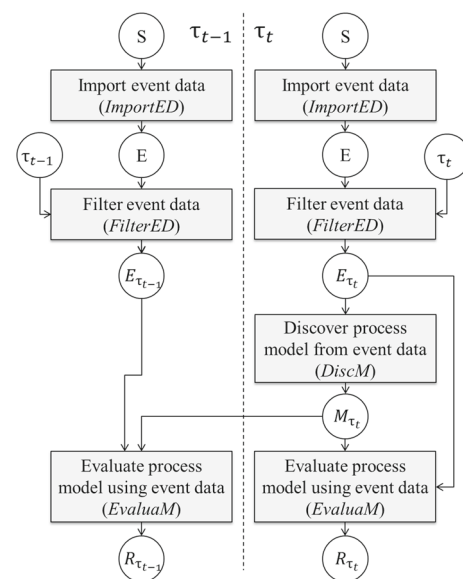
replayed many times using different inputs where no further configuration is required.

There are many examples of this analysis scenario within the process-mining domain. Two representative examples are described next.

### 4.4.1 Periodic benchmarking

Modern organizations make large investments to improve their own processes for better performance in terms of costs, time, or quality. In order to measure these improvements, organizations have to evaluate their performance periodically. This requires them to evaluate performance of the new time-period and compare it with the previous periods. Performance can improve or degrade in different time-periods. Obviously, the returned results require human judgments and, hence, cannot be fully automated by the scientific workflow.

Figure 15 shows an example of this analysis scenario using different process-mining building blocks. Let us assume that we want to compare period  $\tau_k$  with period  $\tau_{k-1}$ . For period  $\tau_k$ , the entire event data are loaded and, then, filtered so as to only keep portion  $E_{\tau_k}$  that refers to the period  $\tau_k$ . Using portion  $E_{\tau_k}$ , a process model  $M_{\tau_k}$  is discovered. For period  $\tau_{k-1}$ , the entire event data are loaded and, then, filtered so as to only keep the portion  $E_{\tau_{k-1}}$  that refers to the period  $\tau_{k-1}$ . Finally, an evaluation is computed about the conformance between model  $M_{\tau_k}$  and event-data portion  $E_{\tau_k}$  and between  $M_{\tau_k}$  and  $E_{\tau_{k-1}}$ . Each evaluation will return valuable results, which are compared to find significant changes.



**Fig. 15** Periodic performance benchmark: process mining workflow for comparing the performance of the process in two different time-periods ( $t$  and  $t - 1$ )



#### 4.4.2 Report generation over collections of data sets

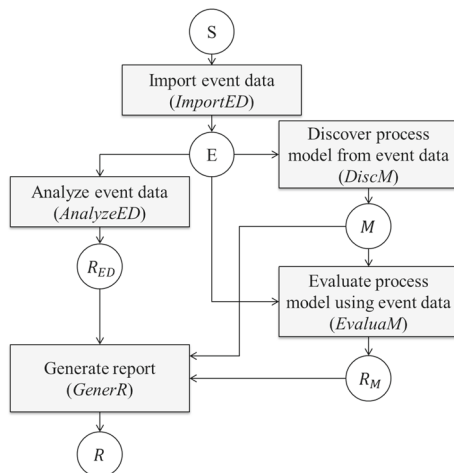
Scientific workflows are useful when generating several reports for different portions of event data, e.g., different groups of patients or customers. Since the steps are the same and the only difference is concerned with using different portions of events, this can be easily automated, even when dozens of subsets need to be taken into consideration.

From this, it follows that this scenario shares common points with large-scale experiments. However, some differences exist. The report-generation scenario is characterized by a stable workflow with a defined set of parameters, whereas in the large-scale experiments scenario, parameters may change significantly in the different iterations. In addition, the input elements used in report-generation scenarios are similar and comparable event data sets. This can be explained by the desire that reports should have the same structure. In case of large-scale experiments, event data sets may be heterogenous. It is thus worthwhile repeating the experiments using diverse and dissimilar event data sets as input.

Figure 16 illustrates a potential scientific workflow to generate reports that contain process-mining results. For the sake of explanation, the process mining workflow is kept simple. The report is assumed to contain three objects: the results  $R_{ED}$  of the analysis of the input event data, the discovered process model  $M$ , and the results  $R_M$  of the evaluation of such a model against the input event data. Process-mining building block *Generate report* takes these three objects as input and combines them into a reporting document  $R$ .

## 5 Implementation

Our framework to support process mining workflows shown in Fig. 1 is supported by *RapidProM*. *RapidProM* was imple-



**Fig. 16** Report generation workflow

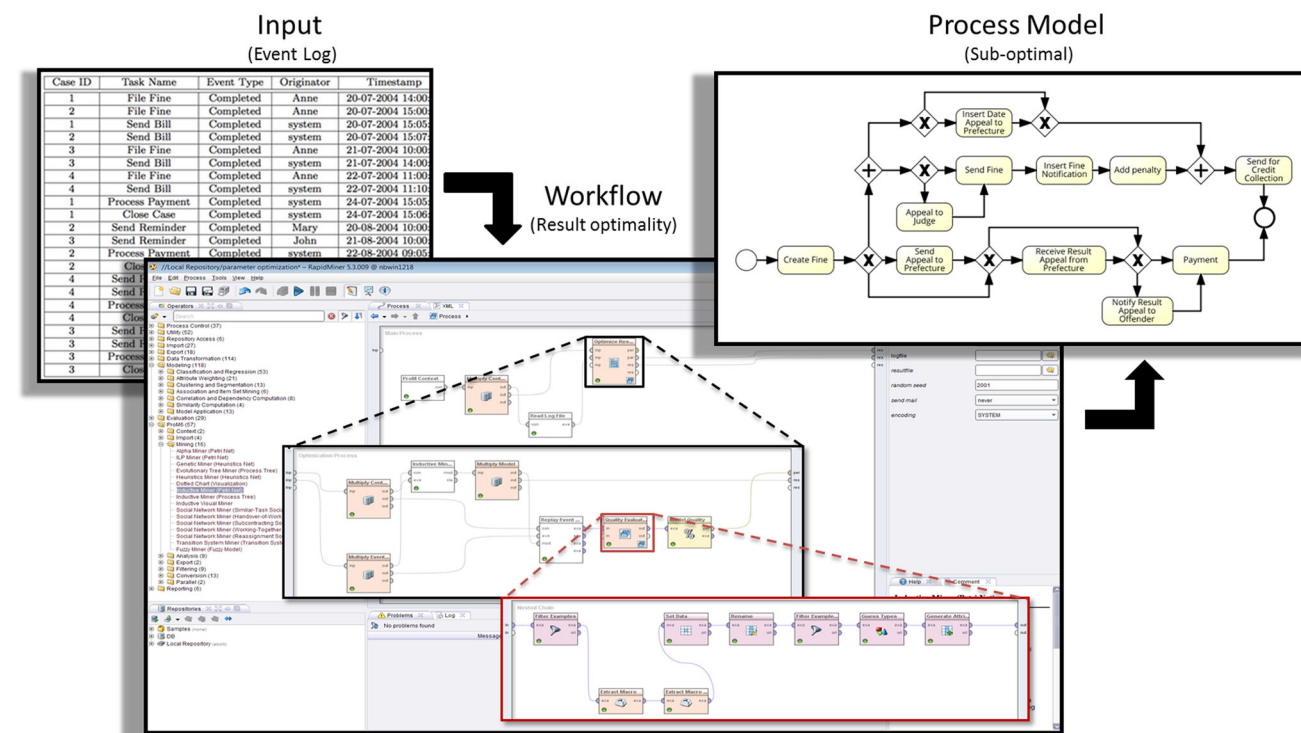
mented using ProM [58] and RapidMiner [22]. The building blocks defined in Sect. 3 have been implemented in RapidProM as *Operators*. Most of the building blocks have been realized using RapidMiner-specific wrappers of plug-ins of the ProM Framework [58]. ProM is a framework that allows researchers to implement process mining algorithms in a standardized environment, which provides a number of facilities to support programmers. Nowadays, it has become the *de-facto* standard for process mining. ProM can be freely downloaded from <http://www.promtools.org>. The extension of RapidMiner to provide process-mining blocks for scientific workflows using ProM is also freely available. At the time of writing, RapidProM provides 37 process mining operators, including several process-discovery algorithms and filters as well as importers and exporters from/to different process-modeling notations. The operators are defined as atomic steps; however, they can be composed into (sub) processes natively in RapidMiner. A (sub) process is the equivalent of a collapsed group of operators, but it can also be executed as an atomic block itself. This is allowed by RapidMiner's native concurrency management, which separates input from output object representations (i.e., a modified input does not affect any other parallel operators that use the same input).

The first version of RapidProM was presented during the BPM 2014 demo session [34]. This initial version successfully implemented basic process-mining functionalities and has been downloaded 4020 times since its release in July 2014 until April 2015 (on average, over 400 monthly downloads). However, process mining is a relatively new discipline, which is developing and evolving rapidly. Therefore, various changes and extensions were needed to keep up with the state-of-the-art. The new version incorporates implementations of various new algorithms, which did not exist in the first version.

The RapidProM extension is hosted both at <http://www.rapidprom.org> and in the RapidProM extension manager server, which can be directly accessed through the RapidMiner Marketplace. After installation, the RapidProM operators are available for use in any RapidMiner workflow, which allows to combine process-mining with other data-mining techniques. Figure 17 shows an example of a process-mining scientific workflow implemented using RapidProM operators. Many of these operators implement a process-mining building block. The process mining workflow shown in Fig. 17 is used in Sect. 6.1 to obtain a sub-optimal process model from event data.

Readers are referred to <http://www.rapidprom.org> for detailed installation, setup and troubleshooting instructions.

Table 1 shows the ProM import plugins implemented in RapidProM Version 2. These five operators are complemented with RapidMiner native operators to export visual results and data tables, in a way that most final results of



**Fig. 17** Example of a process-mining workflow in RapidMiner through the RapidProM extension: the workflow transforms Event data (Input) into a Sub-optimal Process Model (Output)

**Table 1** Import/export operators

Building block	Operator name	Operator description
<i>ImportED</i>	Read Log (path)	Imports an event log from a specified path
<i>ImportED</i>	Read Log (file)	Takes a file object (usually obtained from a “loop files” operator) and transforms it to an Event Log
<i>ImportM</i>	Read PNML	Imports a Petri Net in a <i>Petri Net Modeling Language</i> (PNML) format from a specified path
–	Export Event Log	Exports an Event Log in different formats
–	Export PNML	Exports a Petri Net in PNML format

process mining workflows can be exported and saved outside RapidMiner.

Table 2 shows a list of ProM Discovery plugins implemented in RapidProM as Discovery Operators. These nine operators (usually referred to as *miners*) are the most commonly used discovery techniques for process mining. These discovery operators produce different models using different techniques and parameters to fine-tune the resulting model.

Table 3 shows a list of ProM visualization plugins implemented in RapidProM as visualization operators. These four

visualization plugins are accompanied by renderers that allow one to inspect both intermediate and final results during and after the execution of process mining workflows.

Table 4 shows a list of ProM conversion plugins implemented in RapidProM as conversion operators. These four conversion plugins are intended for converting models into other model formats. This way we improve the chances that a produced model can be used by other operators. For example, if a heuristics net is discovered from an Event Log using the Heuristics Miner, then the Replay Log on Petri Net (Conformance) operator cannot be executed unless a conversion to Petri Net is performed (which is supported).

Table 5 shows a list of log processing operators implemented in RapidProM. Some of these eight operators use ProM functionalities to perform their tasks, but others were developed specifically for RapidProM, as the ProM framework generally does not use flat data tables to represent event data. These operators are used to modify an event log by adding attributes, events, or converting it to data tables, and vice versa.

Table 6 shows a list of ProM plugins implemented in RapidProM as analysis operators.

## 6 Evaluation

This section shows a number of instantiations of scientific workflows in RapidProM, highlighting the benefits of using

**Table 2** Discovery operators

Building block	Operator name	Operator description
<i>DiscM</i>	Alpha Miner [53]	Discovers a Petri Net. Fast but results are not always reliable because of overfitting issues
<i>DiscM</i>	ILP Miner [54]	Discovers a Petri Net by solving <i>Integer Linear Programming</i> (ILP) problems. Result have perfect fitness but generally poor precision. Slow on large Logs
<i>DiscM</i>	Genetic Miner [43]	Discovers a Heuristics Net using genetic algorithms. Depending on the parameter settings it can be slow or fast
<i>DiscM</i>	Evolutionary Tree Miner [11]	Discovers a Process Tree using a guided genetic algorithms based on model quality dimensions. Guarantees soundness but cannot represent all possible behavior due to its block-structured nature
<i>DiscM</i>	Heuristics Miner [62]	Discovers a Heuristics Net using a probabilistic approach. Good when dealing with noise. Fast
<i>DiscM</i>	Inductive Miner [29]	Discovers a Process Tree or Petri Net. Good when dealing with infrequent behavior and large Logs. Soundness is guaranteed
<i>DiscM</i>	Social Network Miner [50]	Discovers a Social Network from the Event Log resources. Different Social Networks can be obtained: similar task, handover of work, etc.
<i>DiscM</i>	Transition System Miner [44]	Discovers a Transition System using parameters to simplify the space-state exploration
<i>DiscM</i>	Fuzzy Miner [18]	Discovers a Fuzzy Model. Good when dealing with unstructured behavior. Fast

**Table 3** Visualization operators

Building block	Operator name	Operator description
<i>AnalyzeED</i>	Dotted Chart [37]	Shows the temporal distribution of events within traces
<i>EnrichM</i>	Inductive Visual Miner [30]	Process exploration tool that shows an annotated interactive model for quick exploration of a Log
<i>EnrichM</i>	Animate Log in Fuzzy Instance [17]	Shows an animated replay of a Log projected over a Fuzzy Instance
<i>EnrichM</i>	PomPom	Petri Net visualizer that emphasizes those parts of the process that correspond to high-frequent events in a given Log

**Table 4** Conversion operators

Building block	Operator name	Operator description
–	Reachability Graph to Petri Net	Converts a Reachability Graph into a Petri Net
–	Petri Net to Reachability Graph	Converts a Petri Net into a Reachability Graph
–	Heuristics Net to Petri Net	Converts a Heuristics Net into a Petri Net
–	Process Tree to Petri Net	Converts a Process Tree into a Petri Net

**Table 5** Log processing operators

Building block	Operator name	Operator description
<i>Added</i>	Add Table Column to Event Log	Adds a single Data Table column as trace attribute to a given Event Log
<i>Added</i>	Add Trace Attributes to Event Log	Adds all columns of a Data Table (except case id) as trace attributes to a given Event Log
<i>Added</i>	Add Event Attributes to Event Log	Adds all columns of a Data Table (except case id and event id) as event attributes to a given Event Log
<i>Added</i>	Add Events to Event Log	Adds Events to a given Event Log from selected columns on a Data Table
<i>MergeED</i>	Merge Event Logs	Merges two Event Logs
<i>Added</i>	Add Artificial Start and End Event	Adds an artificial Start Event to the beginning, and an artificial End Event to the ending of each trace
–	Event Log to ExampleSet	Converts an Event Log into a Data Table (ExampleSet)
–	ExampleSet to Event Log	Converts a Data Table (ExampleSet) into an Event Log

**Table 6** Analysis operators

Building block	Operator name	Operator description
<i>AnalyzeM</i>	WOFLAN [60]	Analyzes the soundness of a Petri Net using the Woflan software ( <a href="http://www.swmath.org/software/7028">www.swmath.org/software/7028</a> )
–	Select fuzzy instance	Selects the best fuzzy instance from a Fuzzy Model
<i>RepairM</i>	Repair Model [15]	Replays an Event Log in a Petri Net and repairs this net to improve fitness
<i>RepairM</i>	Reduce Silent Transitions	Reduces a Petri Net by removing invisible transitions (and places) that are not used
<i>AnalyzeED</i>	Feature Prediction [13]	Produces predictions of business process features using decision trees
<i>EnrichM</i>	Replay Log on Petri Net (Performance) [48]	Replays a Log on a Petri Net and generates performance metrics such as throughput time, waiting time, etc.
<i>EnrichM</i>	Replay Log on Petri Net (Conformance) [48]	Replays a Log on a Petri Net and generates conformance metrics such as fitness

scientific workflows for process mining. They are specific examples of the analysis scenarios discussed in Sect. 4

### 6.1 Evaluating result optimality

The first experiment is related to the *Result Optimality* analysis scenario described in Sect. 4.1. In this experiment, we implemented a process mining workflow using RapidProM to extract the model that scores higher with respect to the geometric average of precision and replay fitness.<sup>2</sup> The geometric average of replay fitness and precision seems to be better than the arithmetic average since it is necessary to have a strong penalty if one of the criteria is low.

For this experiment, we employed the *Inductive Miner - Infrequent* discovery technique [29] and used different values for the *noise threshold* parameter. This parameter is defined in a range of values between 0 and 1. This parameter allows for filtering out infrequent behavior contained in event data

in order to produce a simpler model: the lower the value is for this parameter (i.e., close to 0), the larger the fraction of behavior observed in the event data that the model allows. To measure fitness and precision, we employ the conformance-checking techniques reported in [1,2]. All techniques are available as part of the RapidProM extension. A summary of the concrete operators used for each building block is presented in Table 7.

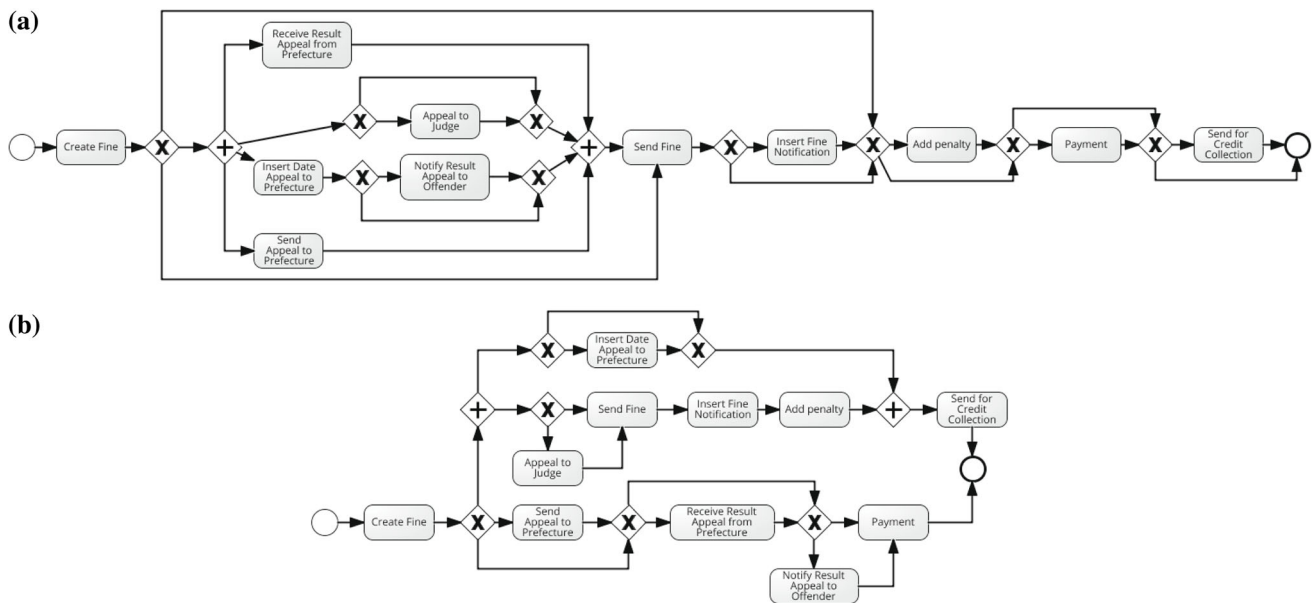
This experiment instantiates the analysis scenario described in Sect. 4.1 and depicted in Fig. 9. The model

**Table 7** Operators used in the result (sub) optimality experiment

Building block	Operator name
<i>ImportED</i>	Read Log (file)
<i>DiscM</i>	Inductive Miner
<i>Evaluam</i>	Replay Log on Petri Net
<i>SelectM</i>	Optimize Parameter (RapidMiner native operator)

<sup>2</sup> The geometric average of  $a$  and  $b$  used in this article is defined as  $2*a*b/(a+b)$  and it is meant to penalize very low scores.





**Fig. 18** Comparison of process models that are mined with the default parameters and with the parameters that maximize the geometric average of replay fitness and precision. The process is concerned with road-traffic fine management and models are represented using the BPMN notation. **a** Model mined using the Inductive Miner with default

value of the noise-threshold parameter, which is 0.2. The geometric average of fitness and precision is 0.708. **b** Model mined using the Inductive Miner with one of the best values of the noise-threshold parameter, which is 0.7. This value was obtained as a result of this experiment. The geometric average of fitness and precision is 0.912

obtained with the default value of the parameter is compared with the model that (almost) maximizes the geometric average of fitness and precision. To obtain this result, we designed a scientific workflow where several models are discovered with different values of the *noise threshold* parameter. Finally, the workflow selects the model with the highest value of the geometric average among those discovered. As input, we used an event-data log that records real-life executions of a process for road-traffic fine management, which is employed by a local-police force in Italy [12]. This event data refer to 150,370 process-instance executions and record the execution of around 560,000 activities.

Figure 18b shows the model generated using the optimal parameters obtained through our scientific workflow, whereas Fig. 18a illustrates the model generated using default parameters.

There are clear differences between the models. For example, in the default model, parallel behavior dominates the beginning of the process. Instead, the “optimal model” presents simpler choices. Another example concerns the final part of the model. In the default model, the latest process activities can be skipped through. However, in the optimal model, this is not possible. The optimal model has a replay fitness and precision of 0.921 and 0.903 respectively, with geometric average 0.912. It scores better than the model obtained through default parameters, where the replay fitness and precision is 1 and 0.548, respectively, with geometric

average 0.708. The optimal model was generated with value 0.7 for the noise threshold parameter.

## 6.2 Evaluating parameter sensitivity

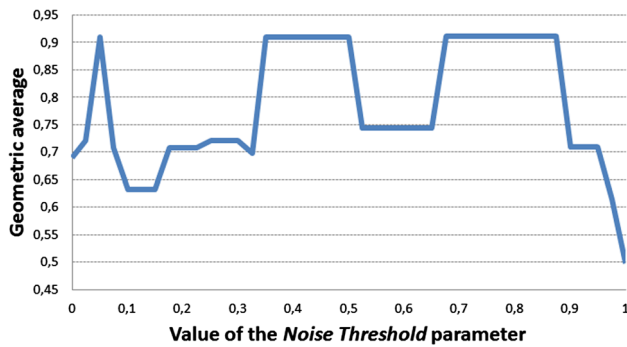
As second experiment illustrating the benefits of using scientific workflows for process mining, we conducted an analysis of the sensitivity of the *noise threshold* parameter of the *Inductive Miner-infrequent*. We used again the event data of the road-traffic fine management process also used in Sect. 6.1. This experiment operationalizes the analysis scenario discussed in Sect. 4.2 and depicted in Fig. 11. In this experiment, we implemented a process mining workflow using RapidProM to explore the effect of this parameter in the final quality of the produced model. In order to do so, we discovered 41 models using different parameter values between 0 and 1 (i.e., a step-size 0.025) and evaluated their quality through the geometric average of replay fitness and precision used before. A summary of the concrete operators used for each building block is presented in Table 8.

Figure 19 shows the results of these evaluations, showing the variation of the geometric average for different values of the *noise threshold* parameter.

By analyzing the graph, the models with higher geometric average are produced when the parameter takes on a value between 0.675 and 0.875. The worst model is obtained when value 1 is assigned to the parameter.

**Table 8** Operators used in the result (sub) optimality experiment

Building block	Operator name
<i>ImportED</i>	Read Log (file)
<i>DiscM</i>	Inductive Miner
<i>CompareM</i>	Replay Event Log on Petri Net, Create chart from value array (RapidMiner native operator)

**Fig. 19** Parameter sensitivity analysis: variation of the geometric average of fitness and precision when varying the value of the *noise threshold* parameter

### 6.3 Performing large scale experiments

As mentioned before, the use of scientific workflows is very beneficial for conducting large-scale experiments with many event logs. When assessing a certain process-mining technique one cannot rely on a single event log to draw conclusions.

For instance, here we want to study how the *noise threshold* parameter influences the quality of the discovered model, in terms of geometric average of fitness and precision. In

Sect. 4.2, the experiment was conducted using a single event log, but RapidProM allows us to do this for any number of event logs using the same operators. To illustrate this, we use 11 real-life event logs and produce the corresponding process models using different parameter settings.

Table 9 shows the results of this evaluation, where each cell shows the geometric average of the replay fitness and the precision of the model obtained using a specific parameter value (column) and event data (row). Every event log used in this experiment is publicly available through the Digital Object Identifiers (DOIs) of the included references. To use some of them for discovery, we had to conduct pre-processing steps in the following cases: the hospital event data set [55] was extremely unstructured. To provide reasonable results and to allow for conformance checking using alignments, we filtered the event log to retain the 80% most frequent behavior before applying the mining algorithm. The same filtering was done for the five CoSeLog event logs [5–9].

The actual results in Table 9 are not very relevant for this paper. It just shows that techniques can be evaluated on a large scale by using scientific workflows.

### 6.4 Automatic report generation

To illustrate the fourth analysis scenario we used event data related to the study behavior and actual performance of students of the Faculty of Mathematics and Computer Science at Eindhoven University of Technology (TU/e). TU/e provides video lectures for many courses to support students who are unable to attend face-to-face lectures for various reasons. The event data record the views of video lectures and the exam attempts of all TU/e courses.

First of all, students generate events when they watch lectures. It is known how long and when they watch a particular

**Table 9** Summary of a few large-scale experimental results: evaluating the geometric average of replay fitness and precision of models discovered with the Inductive Miner using different values of the noise threshold parameter (columns) and different real-life sets of event data (rows)

Event data	<i>nt</i> = 0	<i>nt</i> = 0.1	<i>nt</i> = 0.2	<i>nt</i> = 0.3	<i>nt</i> = 0.4	<i>nt</i> = 0.5	<i>nt</i> = 0.6	<i>nt</i> = 0.7	<i>nt</i> = 0.8	<i>nt</i> = 0.9	<i>nt</i> = 1	Average
BPI2012 [56]	0.314	0.730	0.430	0.450	0.508	0.474	0.675	0.683	0.674	0.679	0.644	0.569
BPI2013 [39]	0.847	0.826	0.778	0.863	0.458	0.458	0.458	0.458	0.458	0.458	0.453	0.592
BPI2014 [57]	0.566	0.720	0.708	0.613	0.616	0.654	0.626	0.414	0.530	0.527	0.490	0.588
Hospital [55]	0.153	0.111	0.546	0.473	0.338	0.172	0.280	0.342	0.392	0.515	0.517	0.349
Road Fines [12]	0.689	0.633	0.708	0.721	0.909	0.909	0.744	0.912	0.912	0.710	0.498	0.758
CoSeLoG 1 [5]	0.143	0.366	0.389	0.576	0.687	0.710	0.737	0.668	0.673	0.649	0.594	0.563
CoSeLoG 2 [6]	0.095	0.191	0.146	0.233	0.127	0.167	0.250	0.177	0.218	0.180	0.362	0.195
CoSeLoG 3 [7]	0.182	0.352	0.573	0.640	0.170	0.209	0.628	0.632	0.585	0.732	0.657	0.487
CoSeLoG 4 [8]	0.190	0.448	0.488	0.640	0.623	0.163	0.553	0.621	0.546	0.518	0.670	0.496
CoSeLoG 5 [9]	0.160	0.199	0.445	0.517	0.522	0.628	0.634	0.145	0.246	0.222	0.602	0.393
CoSeLoG R. [10]	0.520	0.860	0.838	0.869	0.859	0.377	0.868	0.868	0.883	0.861	0.656	0.769
Average	0.350	0.494	0.549	0.599	0.528	0.447	0.586	0.538	0.556	0.550	0.558	

We use *nt* to indicate the value of the *noise threshold* parameter of application of the algorithm

lecture of a particular course. These data can be preprocessed so that low-level events are collapsed into lecture views. Second, students generate events when they take exams and the result is added to the event.

For each course, we have generated a report that includes the results of the application of various data-mining and process-mining techniques. This generation is automatic in the sense that the scientific workflow takes a list of courses as input and produces as many reports as the number of courses in the list.

The report contains three sections: course information, core statistics and advanced analysis.

Figure 20 shows a small part of the report generated for the course on Business Information Systems (2II05). In the first section, the report provides information about the course, the bachelor or master programs which it belongs to, as well as the information about the overall number of views of the course's video lectures. In the second section (only a small fragment is shown), some basic distributions are calculated. For example, statistics are reported about the division per gender, nationality, and final grade. The third

section is devoted to process-mining results. The results of applying conformance checking using the event data and the ideal process model where a student watches every video lecture and in the right order, namely he/she watches the  $i^{th}$  video lecture only after watching the  $(i - 1)^{th}$  video lecture. As expected, the results show a positive correlation between higher grades and higher compliance with the normative process just mentioned: the more a student watches all video lectures in the right order, the higher the corresponding grade will be. In addition to showing the conformance information, the report always embeds a dotted chart. The dotted chart is similar to a Gantt chart (see building block *AnalyzeED*). The dotted chart shows the distribution of events for the different students over time. This way one can see the patterns and frequency with which students watch video lectures.

Note that reports like the one shown in Fig. 20 are very informative for both professors and students. By using RapidProM we are able to automatically generate reports for all courses (after data conversion and modeling the desired process mining workflow).

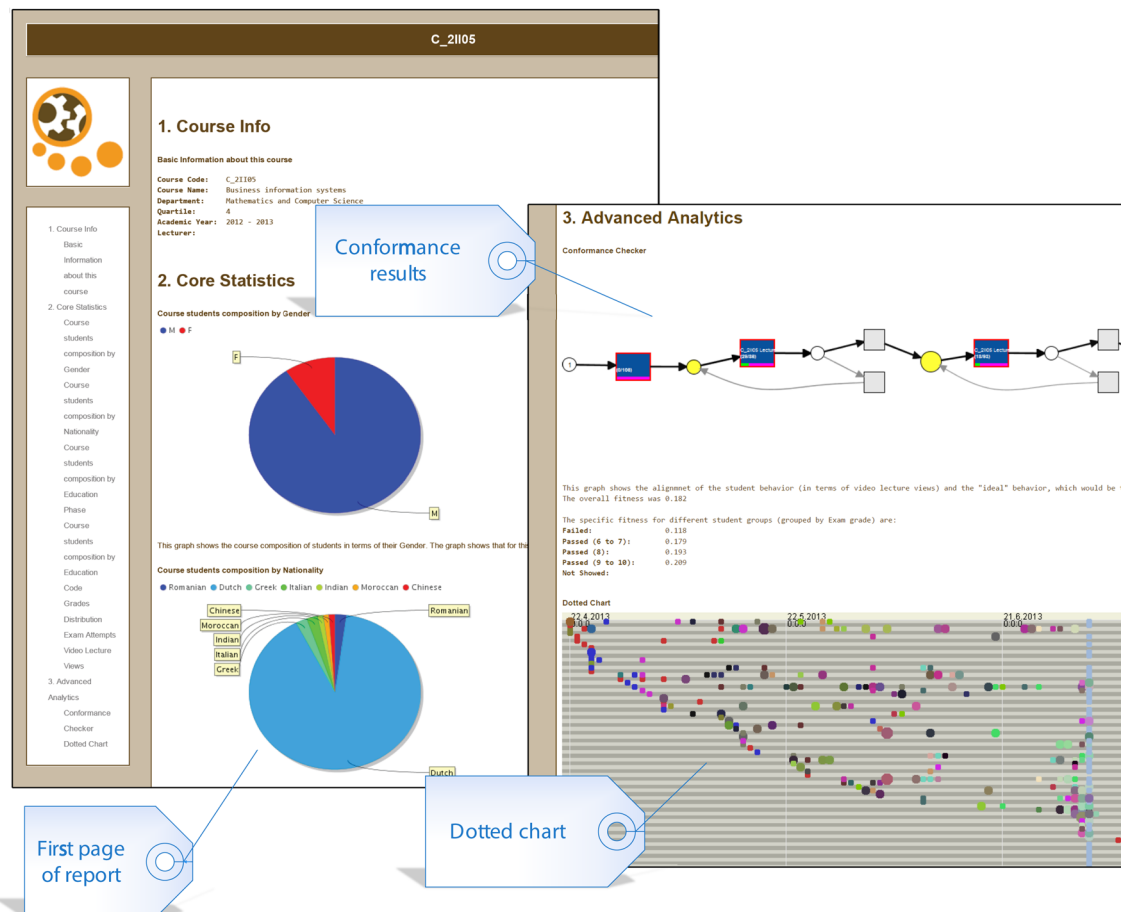


Fig. 20 Fragments of the automatically generated report using RapidProM

## 7 Conclusions

This paper presented a framework for supporting the design and execution of process mining workflows. As argued, scientific workflow systems are not tailored towards the analysis of processes based on models and logs. Tools like RapidMiner and KNIME can model analysis workflows but do not provide any process mining capabilities. The focus of these tools is mostly on traditional data mining and reporting capabilities that tend to use tabular data. Also more classical Scientific Workflow Management (SWFM) systems like Kepler and Taverna do not provide dedicated support for artifacts like process models and event logs. Process mining tools like ProM, Disco, Perceptive, Celonis, QPR, etc. do not provide any workflow support. The inability to model and execute process mining workflows was the primary motivation for developing the framework presented in this paper.

We proposed generic *process mining building blocks* grouped into six categories. These can be chained together to create process mining workflows. We identified four broader *analysis scenarios* and provided conceptual workflows for these. The whole approach is supported using *RapidProM* which is based on ProM and RapidMiner. RapidProM has been tested in various situations and in this paper we demonstrated this using concrete instances of the four analysis scenarios. RapidProM is freely available via <http://www.rapidprom.org> and the RapidMiner Market place.

Future work aims at extending the set of process mining building blocks and evaluating RapidProM in various case studies. We continue to apply RapidProM in all four areas described. Moreover, we would like to make standard workflows available via infrastructures like myExperiment and OpenML. We are also interested a further cross-fertilizations between process mining and other analysis techniques available in tools like RapidMiner and KNIME (text mining, clustering, predictive analytics, etc.).

**Acknowledgments** The authors thank Ronny Mans for his seminal work on the initial version of RapidProM.

**Open Access** This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

## References

- Adriansyah, A., Munoz-Gama, J., Carmona, J., van Dongen, B.F., van der Aalst, W.M.P.: Alignment based precision checking. In: La Rosa, M., Pnina, S. (eds.) *Business Process Management Workshops*. Lecture Notes in Business Information Processing, vol. 132, pp. 137–149. Springer, Heidelberg (2013)
- Adriansyah, A., Munoz-Gama, J., Carmona, J., van Dongen, B.F., van der Aalst, W.M.P.: Measuring precision of modeled behavior. *Inf. Syst. e-bus Manag.* **13**(1), 37–67 (2015)
- Barga, R., Gannon, D.: Scientific versus business workflows. In: Taylor, I.J., Deelman, E., Gannon, D.B., Shields, M. (eds.) *Workflows for e-Science*, pp. 9–16. Springer, Berlin (2007)
- Berthold, M.R., Cebron, N., Dill, F., Gabriel, T.R., Koetter, T., Meinl, T., Ohl, P., Sieb, C., Thiel, K., Wiswedel, B.: *Knime: The konstanz information miner*. In: Preisach, C., Burkhardt, H., Schmidt-Thieme, L., Decker, R. (eds.) *Data Analysis, Machine Learning and Applications, Studies in Classification, Data Analysis, and Knowledge Organization*, pp. 319–326. Springer, Berlin (2008)
- Buijs, J.C.A.M.: Environmental permit application process (wabo), coselog project, municipality 1, (2014). doi:[10.4121/uuid:c45dcbe9-557b-43ca-b6d0-10561e13dcb5](https://doi.org/10.4121/uuid:c45dcbe9-557b-43ca-b6d0-10561e13dcb5)
- Buijs, J.C.A.M.: Environmental permit application process (wabo), coselog project, municipality 2, (2014). doi:[10.4121/uuid:34b4f6f4-dbe0-4857-bf75-5b9e1138eb87](https://doi.org/10.4121/uuid:34b4f6f4-dbe0-4857-bf75-5b9e1138eb87)
- Buijs, J.C.A.M.: Environmental permit application process (wabo), coselog project, municipality 3, (2014). doi:[10.4121/uuid:a8ed945d-2ad8-480e-8348-cf7f06c933b3](https://doi.org/10.4121/uuid:a8ed945d-2ad8-480e-8348-cf7f06c933b3)
- Buijs, J.C.A.M.: Environmental permit application process (wabo), coselog project, municipality 4, (2014). doi:[10.4121/uuid:e8c3a53d-5301-4afb-9bcd-38e74171ca32](https://doi.org/10.4121/uuid:e8c3a53d-5301-4afb-9bcd-38e74171ca32)
- Buijs, J.C.A.M.: Environmental permit application process (wabo), coselog project, municipality 5, (2014). doi:[10.4121/uuid:c399c768-d995-4086-adda-c0bc72ad02bc](https://doi.org/10.4121/uuid:c399c768-d995-4086-adda-c0bc72ad02bc)
- Buijs, J.C.A.M.: Receipt phase of an environmental permit application process (wabo), coselog project. (2014). doi:[10.4121/uuid:a07386a5-7be3-4367-9535-70bc9e77dbe6](https://doi.org/10.4121/uuid:a07386a5-7be3-4367-9535-70bc9e77dbe6)
- Buijs, J.C.A.M., van Dongen, B.F., van der Aalst, W.M.P.: On the role of fitness, precision, generalization and simplicity in process discovery. In: Meersman, R., Panetto, H., Dillon, T., Rinderle-Ma, S., Dadam, P., Zhou, X., Pearson, S., Ferscha, A., Bergamaschi, S., Cruz, I.F. (eds.) *On the Move to Meaningful Internet Systems (OTM)*, volume 7565 of *Lecture Notes in Computer Science*, pp. 305–322. Springer, Berlin (2012)
- de Leoni, M., Mannhardt, F.: Road traffic fine management process. (2015). doi:[10.4121/uuid:270fd440-1057-4fb9-89a9-b699b47990f5](https://doi.org/10.4121/uuid:270fd440-1057-4fb9-89a9-b699b47990f5)
- de Leoni, M., van der Aalst, W.M.P., Dees, M.: A general framework for correlating business process characteristics. In: Sadiq, S., Soffer, P., Vlzer, H. (eds.) *Business Process Management*, volume 8659 of *Lecture Notes in Computer Science*, pp. 250–266. Springer, New York (2014)
- Diamantini, C., Potena, D., Storti, E.: Mining usage patterns from a repository of scientific workflows. In: *Proceedings of the 27th Annual ACM Symposium on Applied Computing, SAC'12*, p. 152–157, New York, NY, 2012. ACM (2012)
- Fahland, D., van der Aalst, W.M.P.: Model repair aligning process models to reality. *Inf Syst* **47**, 220–243 (2015)
- Garijo, D., Alper, P., Belhajjame, K., Corcho, Ó., Gil, Y., Goble, C.A.: Common motifs in scientific workflows: an empirical analysis. *Future Gener. Comput. Syst.* **36**, 338–351 (2014)
- Günther, C.W.: *Process Mining in Flexible Environments*, PhD thesis. Technische Universiteit Eindhoven, Eindhoven, The Netherlands (2009)
- Günther, C.W., van der Aalst, W.M.P.: Fuzzy mining adaptive process simplification based on multi-perspective metrics. In: Alonso, G., Dadam, P., Rosemann, M. (eds.) *Business Process Management*, volume 4714 of *Lecture Notes in Computer Science*, pp. 328–343. Springer, Berlin (2007)
- Goble, C.A., Bhagat, J., Aleksejevs, S., Cruickshank, D., Michaelides, D., Newman, D., Borkum, M., Bechhofer, S., Roos, M., Li, P., De Roure, D.: myExperiment: a repository and social



- network for the sharing of bioinformatics workflows. *Nucl. Acids Res.* **38**(suppl 2), W677–W682 (2010)
20. Goecks, J., Nekrutenko, A., Taylor, J., The Galaxy Team.: Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences. *Genome Biol* **11**(8), R86 (2010)
  21. Hand, D.J., Smyth, P., Mannila, H.: *Principles of Data Mining*. MIT Press, Cambridge, MA (2001)
  22. Hofmann, M., Klinkenberg, R.: *RapidMiner: Data Mining Use Cases and Business Analytics Applications*. Chapman and Hall/CRC, Florida (2013)
  23. Hull, D., Wolstencroft, K., Stevens, R., Goble, C.A., Pocock, M.R., Li, P., Oinn, T.: Taverna: a tool for building and running workflows of services. *Nucl. Acids Res.* **34**, 729–732 (2006)
  24. IEEE Task Force on Process Mining. Process mining case studies. [http://www.win.tue.nl/ieeetfpm/doku.php?id=shared:process\\_mining\\_case\\_studies](http://www.win.tue.nl/ieeetfpm/doku.php?id=shared:process_mining_case_studies) (2013)
  25. Ihaka, R., Gentleman, R.: R: a language for data analysis and graphics. *J. Comput. Graph. Stat.* **5**(3), 299–314 (1996)
  26. Keim, D., Andrienko, G., Fekete, J.-D., Grg, C., Kohlhammer, J., Melancon, G.: Visual analytics: definition, process, and challenges. In: Kerren, A., Stasko, J.T., Fekete, J.-D., North, C. (eds.) *Information Visualization*, volume 4950 of *Lecture Notes in Computer Science*, pp. 154–175. Springer, Berlin (2008)
  27. Kranjc, J., Podpean, V., Lavra, N.: Clowdflows: a cloud based scientific workflow platform. In: Flach, P.A., De Bie, T., Cristianini, N. (eds.) *Machine Learning and Knowledge Discovery in Databases*, volume 7524 of *Lecture Notes in Computer Science*, pp. 816–819. Springer, Berlin (2012)
  28. La Rosa, M., Reijers, H.A., van der Aalst, W.M.P., Dijkman, R.M., Mendling, J., Dumas, M., García-Bañuelos, L.: Apomore: an advanced process model repository. *Expert Syst. Appl.* **38**(6), 7029–7040 (2011)
  29. Leemans, S.J.J., Fahland, D., van der Aalst, W.M.P.: Discovering block-structured process models from event logs—a constructive approach. In: Colom, J.-M., Desel, J. (eds.) *Application and Theory of Petri Nets and Concurrency*, volume 7927 of *Lecture Notes in Computer Science*, pp. 311–329. Springer, Berlin (2013)
  30. Leemans, S.J.J., Fahland, D., van der Aalst, W.M.P.: Exploring processes and deviations. In: Fournier, F., Mendling, J. (eds.) *Business Process Management Workshops*, volume 202 of *Lecture Notes in Business Information Processing*, pp. 304–316. Springer, Heidelberg (2015)
  31. Leymann, F., Roller, D.: *Production Workflow—Concepts and Techniques*. Prentice Hall PTR, Upper Saddle River, NJ (2000)
  32. Littauer, R., Ram, K., Ludäscher, B., Michener, W., Koskela, R.: Trends in use of scientific workflows: insights from a public repository and recommendations for best practice. *IJDC* **7**(2), 92–100 (2012)
  33. Ludäscher, B., Altintas, I., Berkley, C., Higgins, D., Jaeger, E., Jones, M.B., Lee, E.A., Tao, J., Zhao, Y.: Scientific workflow management and the kepler system. *Concurr Comput. Pract. Exp.* **18**(10), 1039–1065 (2006)
  34. Mans, R.S., van der Aalst, W.M.P., Verbeek, H.M.W.: Supporting process mining workflows with RapidProM. In: Limonad, L., Weber, B. (eds.) *Proceedings of the BPM Demo Sessions 2014 Collocated with the 12th International Conference on Business Process Management (BPM)*, volume 1295 of *CEUR Workshop Proceedings*, pp. 56–60. CEUR-WS.org (2014)
  35. Mitchell, T.M.: *Machine Learning*. McGraw Hill Series in Computer Science. McGraw-Hill, New York (1997)
  36. Scheer, A.-W., Nüttgens, M.: Aris architecture and reference models for business management. In: van der Aalst, W.M.P., Desel, J., Oberweis, A. (eds.) *Business Process Management*, volume 1806 of *Lecture Notes in Computer Science*, pp. 376–389. Springer, Berlin (2000)
  37. Song, M., van der Aalst, W.M. P.: Supporting process mining by showing events at a glance. In: *Proceedings of the 17th Annual Workshop on Information Technologies and Systems (WITS)*, pp. 139–145 (2007)
  38. Sonntag, M., Karastoyanova, D., Deelman, E.: Bridging the gap between business and scientific workflows: humans in the loop of scientific workflows. In: *Sixth International Conference on e-Science, e-Science 2010*, 7–10 Dec 2010, Brisbane, QLD, pp. 206–213 (2010)
  39. Steeman, W.: Bpi challenge 2013. (2013). doi:[10.4121/500573e6-acc-4b0c-9576-aa5468b10cee](https://doi.org/10.4121/500573e6-acc-4b0c-9576-aa5468b10cee)
  40. Steffen, B., Margaria, T., Nagel, R., Joerges, S., Kubczak, C.: Model-driven development with the jABC. In: Bin, E., Ziv, A., Ur, S. (eds.) *Hardware and Software, Verification and Testing*, volume 4383 of *Lecture Notes in Computer Science*, pp. 92–108. Springer, Berlin (2007)
  41. Taylor, I.J., Deelman, E., Gannon, D.B., Shields, M.: *Workflows for e-Science: Scientific Workflows for Grids*. Springer, Berlin (2007)
  42. Turner, K.J., Lambert, P.S.: Workflows for quantitative data analysis in the social sciences. *Int. J. Softw. Tools Technol. Transf.* **17**(3), 321–338 (2015)
  43. van der Aalst, W.M.P., Alves de Medeiros, A.K., Weijters, A.J.M.M.: Genetic process mining. In: Ciardo, G., Darondeau, P. (eds.) *Applications and Theory of Petri Nets 2005*, volume 3536 of *Lecture Notes in Computer Science*, pp. 48–69. Springer, Berlin (2005)
  44. van der Aalst, W.M.P., Schonenberg, M.H., Song, M.: Time prediction based on process mining. *Inf. Syst.* **36**(2), 450–475 (2011). Special Issue: Semantic Integration of Data, Multimedia, and Services
  45. van der Aalst, W.M.P.: *Process Mining: Discovery, Conformance and Enhancement of Business Processes*, 1st edn. Springer, Berlin (2011)
  46. van der Aalst, W.M.P.: A decade of business process management conferences: personal reflections on a developing discipline. In: Barros, A., Gal, A., Kindler, E. (eds.) *Business Process Management*, volume 7481 of *Lecture Notes in Computer Science*, pp. 1–16. Springer, Berlin (2012)
  47. van der Aalst, W.M.P.: Decomposing process mining problems using passages. In: Haddad, S., Pomello, L. (eds.) *Application and Theory of Petri Nets*, volume 7347 of *Lecture Notes in Computer Science*, pp. 72–91. Springer, Berlin (2012)
  48. van der Aalst, W.M.P., Adriansyah, A., van Dongen, B.F.: Replaying history on process models for conformance checking and performance analysis. *Wiley interdisciplinary reviews. Data Min. Knowl. Disc.* **2**(2), 182–192 (2012)
  49. van der Aalst, W.M.P., Dreiling, A., Gottschalk, F., Rosemann, M., Jansen-Vullers, M.H.: Configurable process models as a basis for reference modeling. In: Bussler, C.J., Haller, A. (eds.) *Business Process Management Workshops*, volume 3812 of *Lecture Notes in Computer Science*, pp. 512–518. Springer, Berlin (2006)
  50. van der Aalst, W.M.P., Song, M.: Mining social networks: uncovering interaction patterns in business processes. In: Desel, J., Pernici, B., Weske, M. (eds.) *Business Process Management*, volume 3080 of *Lecture Notes in Computer Science*, pp. 244–260. Springer, Berlin (2004)
  51. van der Aalst, W.M.P., ter Hofstede, A.H.M., Kiepuszewski, B., Barros, A.P.: Workflow patterns. *Distrib. Parallel. Databases* **14**(1), 5–51 (2003)
  52. van der Aalst, W.M.P., van Hee, K.M., ter Hofstede, A.H.M., Sidorova, N., Verbeek, H.M.W., Voorhoeve, M., Wynn, M.T.: Soundness of workflow nets: classification, decidability, and analysis. *Form. Aspect. Comput.* **23**(3), 333–363 (2011)
  53. van der Aalst, W.M.P., Weijters, A.J.M.M., Maruster, L.: Workflow mining: discovering process models from event logs. *IEEE Trans. Knowl. Data Eng.* **16**(9), 1128–1142 (2004)

54. van der Werf, J.M.E.M., van Dongen, B.F., Hurkens, C.A.J., Serebrenik, A.: Process discovery using integer linear programming. In: van Hee, K.M., Valk, R. (eds.) *Applications and Theory of Petri Nets*, volume 5062 of *Lecture Notes in Computer Science*, pp. 368–387. Springer, Berlin (2008)
55. van Dongen, B.F.: Real-life event logs - hospital log. (2011). doi:[10.4121/uuid:d9769f3d-0ab0-4fb8-803b-0d1120ffc54](https://doi.org/10.4121/uuid:d9769f3d-0ab0-4fb8-803b-0d1120ffc54)
56. van Dongen, B.F.: Bpi challenge 2012. (2012). doi:[10.4121/uuid:3926db30-f712-4394-aebc-75976070e91f](https://doi.org/10.4121/uuid:3926db30-f712-4394-aebc-75976070e91f)
57. van Dongen, B.F.: Bpi challenge 2014. (2014). doi:[10.4121/uuid:d5ccb355-ca67-480f-8739-289b9b593aaf](https://doi.org/10.4121/uuid:d5ccb355-ca67-480f-8739-289b9b593aaf)
58. van Dongen, B.F., de Medeiros, A.K.A., Verbeek, H.M.W., Weijters, A.J.M.M., van der Aalst, W.M.P.: The ProMframework: a new era in process mining tool support. In: Ciardo, G., Darondeau, P. (eds.) *Applications and Theory of Petri Nets*, volume 3536 of *Lecture Notes in Computer Science*, pp. 444–454. Springer, Berlin (2005)
59. van der Aalst, W.M.P.: Decomposing Petri nets for process mining: a generic approach. *Distrib. Parallel Databases* **31**(4), 471–507 (2013)
60. Verbeek, H.M.W., Basten, T., van der Aalst, W.M.P.: Diagnosing workflow processes using woflan. *Comput. J.* **44**(4), 246–279 (2001)
61. Wassink, I.H.C., van der Vet, P.E., Wolstencroft, K., Neerincx, P.B.T., Roos, M., Rauwerda, H., Breit, T.M.: Analysing scientific workflows: why workflows not only connect web services. In: 2009 IEEE Congress on Services, Part I, SERVICES I 2009, Los Angeles, CA, July 6–10, pp 314–321 (2009)
62. Weijters, A.J.M.M., van der Aalst, W.M.P.: Rediscovering workflow models from event-based data using Little Thumb. *Integr. Comput. Aided Eng.* **10**(2), 151–162 (2003)
63. Weske, M.: *Business Process Management—Concepts, Languages, Architectures*, 2nd edn. Springer, Heidelberg (2012)
64. Wickert, A., Lamprecht, A.-L.: jABCstats: an extensible process library for the empirical analysis of jABC workflows. In: Margaria, T., Steffen, B. (eds.) *Leveraging Applications of Formal Methods, Verification and Validation. Specialized Techniques and Applications*, volume 8803 of *Lecture Notes in Computer Science*, pp. 449–463. Springer, Berlin (2014)
65. Zeng, R., He, X., Li, J., Liu, Z., van der Aalst, W.M.P.: A method to build and analyze scientific workflows from provenance through process mining. In: 3rd Workshop on the Theory and Practice of Provenance, TaPP'11, Heraklion, Crete, Greece, 20–21 June 2011 (2011)